

## REMARKS

The original specification has been replaced with the attached replacement specification to add headings and amend the title.

In amended Fig. 1 and Fig. 4 the blocks have been labeled with descriptive words or abbreviations.

Claims 1-10 remain in the application.

The Examiner has objected to the Title, the Drawings and the Specification. As stated above, these have all been amended. Therefore, it is respectfully requested that the Examiner withdraw his objections.

Claims 8 and 9 were rejected under 35 U.S.C. 102(b) as being unpatentable over Roucos (Speaker Normalization Algorithms for Very-Low-Rate Speech Coding), Fette (4,707,858) or Taguchi (4,701,955). Claims 2-7 were rejected under 35 U.S.C. 103(a) as being unpatentable over Roucos, Fette or Taguchi. Claim 10 was rejected under 35 U.S.C. 103(a) as being unpatentable over Roucos, Fette or Taguchi and further in view of Zehave (5,581,575). These rejections are respectfully traversed.

Roucos merely teaches an algorithm for progressively adapting vocal references, of a given library, to the voice of a speaker, by comparison with power spectral densities and by use of covariance matrices through a fifteen minutes testing. In Roucos, the differences are used for adapting the vocal references of the library and can not be used in real time to the coded words of the speech as secondary codes, as required in Applicants claim 1.

Fette merely teaches a voice recognition system for identifying a user original voice with respect to words in storage, and specifically the methods and apparatus used to insure the correct recognition or to postpone the decision as to the identity of the speaker (column 1, lines 24-50). In Fette, the differences are not used to form code words for transmission, as stated by the Examiner, but are used to determine if the received voice spectrum is noise spectrum or not (for example, see column 3, lines 41-44 and 54-56). This is a major difference. The so called ADMF (Fig. 2) is described in copending US patent application Ser. No. 309,640 and is not described in this document. As a matter of fact, nothing can be derived but the derivation indicated above, i.e. to determinate if the signal in entry is noise or not, because it is only supplied by ADC 30 (see fig. 2), which is a classical function.

Tagushi merely teaches a variable frame length vocoder in which the input speech signal is divided in inclined and not-inclined time sections, the former being of fixed time length and the latter of variable time by trapezoidal approximation, of the LSP coefficient vector (column 4, lines 40-45). Therefore, the speech is coded only with these new coefficients taken as unique code words, and without any secondary codes derived from any library.

Therefore, none of the three above references alone or in combination teaches the requirements of claim 1, that is to encoding speech by primary and secondary codes to from acoustic units. Claim 1 related to a new method to achieve a personalized restoration by use of compressed information, in real time, consisting in doubling the coding information, this being not obvious at all.

Claim 8 is distinguishable over the prior art for the same reasons as claim 1. The comparator means and the transcoding means, with their function, are not taught by any of the three above references.

Claim 9 is distinguishable for the same reasons as well. Moreover, none of the prior references teaches the decoding means (33) with its correction function (38).

Therefore, withdrawal of these rejections are respectfully requested.

It is believed that all claims are in condition for allowance. Early and favorable action by the Examiner is earnestly solicited.

### AUTHORIZATION

A one month extension of time is included with this Amendment. In the event that an extension of time may be required in addition to that requested in a petition for an extension of time, the Commissioner is requested to grant a petition for that extension of time which is required to make this response timely and is hereby authorized to charge any fee for such an extension of time or credit any overpayment for an extension of time to Deposit Account No. 50-1561.

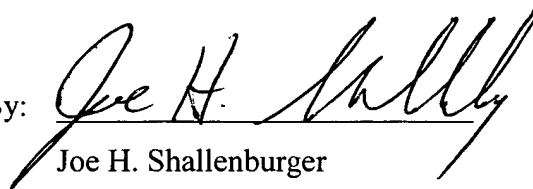
The Commissioner is hereby authorized to charge any additional fees which may be required for this amendment, or credit any overpayment to Deposit Account No. 50-1561.

If the Examiner believes that issues may be resolved by telephone interview, the Examiner is respectfully urged to telephone the undersigned at (212) 801-2146. The undersigned may also be contacted by e-mail at [ecr@gtlaw.com](mailto:ecr@gtlaw.com).

Respectfully submitted,

Dated: November 3, 2003

By:



Joe H. Shallenburger

Registration No. 37,937

TITLE OF THE INVENTION

**Process for encoding speech using primary code words  
and terminals for implementing the process.**

BACKGROUND OF THE INVENTION

The transmission of speech on the switched telephone network STN necessitates a pass-band sufficient for the speech to remain comprehensible. A band ranging from very low frequencies to some kilohertz represents a good compromise between fidelity of restoration and pass-band resources. For this reason, in order to transmit the voice on the STN, of which the inter-centre connections are digital, the voice frequencies are encoded to transform them into a digital signal at the basic rate of the STN, representing the evolution of the amplitude of the voice signal.

However, it is sometimes desirable to reduce the rate of the transmission, for example, in a voice synthesis terminal, of which the message memory must remain of limited size. In the same way, it may be desired to use only one transmission channel of a specific network, with an output well below the afore-mentioned rate of the telephone network.

In such a case, an attempt is made to recognise the various phonemes of a voice sequence instantly. These phonemes are compared with references in a library, which are associated with code words and these phonemes are replaced by the corresponding code words which describe the speech using a much smaller quantity of information. In this way the voice is compressed.

During reception, the called terminal comprises the same library and, by means of voice synthesis, restores analogue signals corresponding to the various code words.

However, such a procedure presents the disadvantage of restoring only a voice which has been standardised by the library and which is thus impersonal, and it is therefore essentially impossible to recognise the correspondent in order to authenticate a voice message. The

inflexions or fluctuations of the voice, which form part of the information just as much as the meaning of the words themselves, are thus not restored.

### BRIEF SUMMARY OF THE INVENTION

The present invention aims to achieve voice encoding which makes it possible both to compress the information and to achieve a personalised restoration.

To this end, the invention firstly relates to a process for encoding speech formed from a sequence of acoustic units, in which the units are compared with library references associated with primary code words, the differences between the units and the references are determined, the differences are encoded by secondary code words and pairs of primary and secondary codes are substituted for the units.

Thus the primary code words will effectively and compactly encode the largest part of the acoustic energy input, while the secondary code words will improve the fidelity of restoration but without greatly increasing the volume of code data since they relate to only a limited energy and a low number of bits makes it possible to encode this marginal energy modulating the primary, standard energy corresponding to the primary code words.

The invention also relates to a terminal for encoding speech signals, comprising means for inputting a sequence of acoustic units and transmitting it to comparator means arranged to compare the acoustic units successively with library references and thus select therein in each instance a specific primary code word of one of the references, the terminal being characterised by the fact that the comparator means are arranged to determine a difference between the input acoustic unit considered and the reference corresponding to the code word selected and to transmit this difference to transcoding means provided to supply, in response, a secondary code word corresponding to memory means arranged to associate the respective primary and secondary code words.

Finally, the invention relates to a terminal for decoding signals, comprising means for receiving signals representing primary code words of references of acoustic units in a library, and decoding means arranged to select certain ones of the references in the library according to the primary code words received and to control a transducer for restoration of speech signals accordingly, the terminal being characterised in that the decoding means are arranged also to

decode secondary, correction code words associated with the primary code words, and to correct the selected voice references accordingly.

Although the process of the invention makes it necessary, all in all, to have an encoding terminal and a corresponding decoding terminal, each of these can be sold separately and the applicant thus intends to claim each one.

In particular, it is advantageous to provide a facsimile machine comprising means for inserting the code words into a facsimile message.



**BRIEF DESCRIPTION OF THE SEVERAL VIEWS OF THE DRAWINGS**

The invention will be better understood with the aid of the following description of a preferred embodiment of the process of the invention, with reference to the attached drawing in which:

- Figure 1 schematically illustrates a transmission terminal and a reception terminal for voice signals for implementation of the process of the invention,
- Figure 2 shows the amplitude  $A$  of a speech signal over time  $t$ ,
- Figure 3 shows the amplitude  $K$  of lines of the spectrum of the speech signal 2 over the frequency  $F$ , and
- Figure 4 is a flow diagram showing the steps of the process.

### DETAILED DESCRIPTION OF THE INVENTION

The transmission terminal referenced 15 in Figure 1, which is in this case in the form of a portable handset for a radio communication network, comprises a microphone 26 for inputting the speech signal of its user, supplying an analogue/digital converter 27 connected at the output to a central microprocessor unit 28 associated with two libraries 11 and 12 of sound sequences or standardised acoustic units such as phonemes. The central unit 28 which encodes the speech controls a transmitter 29, in this case a radio transmitter, of which the transmissions are received by a reception circuit 30 of a speech restoration terminal 35. Figure 2 illustrates the amplitude  $\sim f$  of an acoustic unit according to time  $t$  and Figure 3 illustrates the amplitude  $K$  of the lines of the spectrum corresponding to a given moment.

More precisely, the central unit 28 comprises a comparator 16 to compare the acoustic units received from the converter 27 with library acoustic units. As explained in more detail with regard to Figure 4, the comparator 16 has the function of selecting the library reference which is the most similar to the signal currently analysed and also has the function of specifying this difference, i.e. of providing a deviation value for each of the criteria which serve in the selection. This difference is transcoded by a transcoding circuit 17, in order to condense its expression into the form of a secondary code word which is stored in a memory 18 under the control of the comparator 16. This comparator, which has previously stored the primary code word in the memory 18, addresses and controls the writing therein so that the two code words, primary and secondary, are physically associated, as they are from a logic point of view, i.e. that, for example, chaining is defined between the two memory zones containing them.

The reception terminal 35 comprises a central unit 33 carrying out inverse decoding of the speech to supply a loudspeaker 34. Two memories forming libraries 31 and 32, in this case external to the central unit 33, are connected to this central unit. The reception terminal 35 is in this case a standard terminal for receiving written messages, known as a pager, also arranged to receive voice messages. For the sake of clarity various standard dialling input, dialling transmission and data display circuits have not been shown.

The central unit 33 comprises a circuit 36 for addressing the libraries 31 and 32, the personalised library and primary library respectively, on the basis of code words received from

the reception circuit 30. In response, a buffer circuit 37 receives, from the primary library 32, spectra of primary acoustic units and transmits them to a circuit 38 for modulation or composition of these spectra. The circuit 38 modulates these spectra according to the secondary code word associated with the primary code word read from the primary library 32. The circuit 38 thus combines the information of the primary and secondary code words to restore the speech signal initially captured (26). This combination can, for example, be an addition or multiplication of frequency lines followed by an inverse Fourier transformation or it can also relate directly to signal amplitudes. In this example, each type of restored acoustic unit is stored in the personalised memory 31 in order to use this latter directly if an identical pair of code words, primary and secondary, is later received. In one variation, the memory 31 could contain only modulation values which it would supply to the circuit 38 after addressing by a secondary code word.

The encoding and decoding operations will now be explained in more detail with reference to Figure 4.

In order to encode the voice a speech signal 26 is captured by the microphone 26 during a step 1 and, in this case, it is converted into a digital signal in the converter 27 in a step 2. The speech signal is then compared, in the central unit 28, with a plurality of reference signals from the library 11 in a step 3. The comparison is carried out instantly, in practice in a cyclic manner at high speed with respect to the speed of evolution of the analysed speech signal. This signal can be considered as being a sequence of acoustic units specific to a given language, such as vowels, diphthongs or hiatus, of which a representation has initially been placed in the library 11 and associated with a code word referred to as the primary code word, which is peculiar to each. When the library 11 and the libraries 12 and 32 mentioned hereinunder are being formed, a number of voice inputs from a single speaker or from several are carried out in order to form an average voice reference. However, in order to improve the efficacy of future recognition, a number of references are preferably stored (11, 12) for each acoustic unit in order to form a range of recognition permitting differences between speakers to be accommodated.

Each acoustic unit (Figure 2) corresponds to a particular evolution of the amplitude A or energy of the speech signal and has a duration able to vary according to the talking speed of the person speaking.

The step 3 thus consists of comparing the evolution of the amplitudes of the reference signals to that of the captured signal. In order to overcome variations in the talking speed, it is possible, for example, to consider only the succession of significant modulations of amplitude (variation of energy exceeding a threshold) without associating a time value notion therewith.

In Figure 2 the vertical arrows, of which there are eight, represent the amplitude of the extremes and thus form a signature, assumed, in this case, to represent a specific acoustic unit. If, leaving the time domain of Figure 2, the frequency domain is now considered, the Fourier transform of the momentary amplitude A of the signal at any point on the curve of Figure 2 can be shown by the spectrum of frequency lines of Figure 3. In practice, it is considered that the voice energy is essentially limited to three frequency bands located respectively around 0.1 kHz and two bands between about 1 and 3 kHz as well as 5 and 7 kHz respectively.

For this reason, if the curve of the amplitudes of Figure 2 is traced over time t, the amplitude K of each line of Figure 3 will be modulated according to the evolution of the amplitude A of the speech signal.

Thus, if the succession of the spectra of Figure 3 is stored, it is possible to restore the succession of the amplitudes A of the original signal by an inverse Fourier transform.

In order to limit the number of spectra to be processed, it is possible to effect only cyclic sampling, sufficiently close not to lose information. It is also possible to be limited to a narrow set of spectra of the extremes of amplitude which are shown by the eight arrows in Figure 2. If it is desired to limit the number of spectra further, it is even possible to retain only a single spectrum representing the average of all the spectra over the considered period of time of the acoustic unit or the average of the spectra of the extremes.

As explained hereinunder in more detail, the average spectrum, or the spectra of the speech signal captured will be compared with one or more counterpart spectra of reference speech signals in the library in order, on the one hand, to select the reference speech signal (acoustic unit) which is the most similar to the captured signal and, on the other hand, to produce a signal representing the difference between the spectrum or spectra of this latter and the spectrum or spectra of the selected reference signal. The difference signal is formed into a code word, referred to as a secondary code word, and is associated with the primary code word of the selected reference signal (recognised acoustic unit) and thus constitutes an additional item of information for modulation or correction of the standardised analogue signal which will be restored from the primary code word considered.

The primary code words of the acoustic units, successively selected as the voice sequence progresses, are stored in a step 4 in order to form a message encoded according to the standard of the library 11.

Furthermore, in a step 5 some of the acoustic units captured and recognised are processed more by analysing their frequency spectrum in detail, in this case in the frequency domain by an inverse Fourier transform, as explained above, step 6.

In a step 7, the spectrum of lines  $j$  of the acoustic unit of identity  $I$  concerned or the spectra representing its evolution over time  $t$  are compared with the spectrum or spectra of the acoustic unit selected in the library 11 and which is or are contained in the associated library 12. In this way, for the or each spectrum a series of weighting coefficients  $C_{ijt}$  ( $i$  = identity of phoneme,  $j$  = frequency rank of the line,  $t$  = time rank) is established, each indicating the amplitude or relative energy of each line  $j$  with respect to its counterpart in the library 12. In other words, these coefficients also represent, albeit indirectly, the relative difference  $(1 - C_{ijt})$  between the recognised acoustic unit and the corresponding reference in the library). The lines in each of the three bands actually correspond to a row of mini bands of adjacent frequencies, in which the voice energy is detected. The analysis in the frequency domain, which is selected in this case, thus provides a more detailed item of information than in the case of an analysis in the time domain of Figure 2 where only the momentary amplitude  $A$  is available.

Thus, in the case of Figures 2 and 3, the series above comprises twelve coefficients representing the twelve lines shown, so that a table of eight such series represents the acoustic unit through the eight extremes shown. Apart from the reduction of the table to a single series, it is possible to make provision to retain only a single average weighting coefficient for each of the three bands. If each coefficient is encoded on just 4 bits, the error will not exceed about 3%, which is amply sufficient to restore a voice timbre, especially since the correction signal represents little energy with respect to the normed signal which it corrects, so that the error viewed as a total is low.

It is thus possible in this case to associate with the primary code word of the acoustic unit selected, of the order of about a hundred bits ( $12 \times 8$ ) if each extreme is retained, or only 12 bits ( $4 \times 3$ ) for the three bands. As the voice timbre is particularly provided by the high frequency of the third band, it is even possible to transmit only the secondary, correction code word relating thereto.

In a step 8, the signal representing the difference in the spectra is transformed into a secondary code word representing the table or the series mentioned above. When the captured speech sequence ends, the primary code words of step 4 and the secondary code words of step 8 are associated one to one (step 9) then transmitted on a transmission network such as, for example, the switched telephone network or, in this case, a radio messaging network (step 10).

The called terminal 35 receives the message in a step 21 and, in a step 22, a primary library file 32 similar to the file of spectra 12, is read by the circuit 36 in order to extract therefrom the primary standardised spectra according to the primary code words. In a step 23 the secondary code words serve to modulate (38) the amplitudes or energies of the standardised lines read in the primary library 32 in order, in this way, to constitute the personalised library 31 of acoustic units, i.e. comprising, in particular, the timbre of the captured voice. The acoustic units of the personalised library 31 are represented in digital form in the time domain after previous transformation by an inverse Fourier transform.

In a step 24, the primary code words received are read successively in order to restore the captured speech signal via the loudspeaker 34 (step 25). For this purpose the primary code words read the personalised library 31 which thus corresponds to the library 11 but which has been personalised by the features in the spectrum of the captured voice.

As stated above, the formation of the library 31 is optional and aims to store a correction for each primary code word, which avoids having to repeat the sending of the secondary code word when the same primary code word is transmitted several times, if, on the other hand, a secondary code word is transmitted systematically, this code word can evolve to follow the possible evolutions in the timbre. In this case the restored voice is both personalised and the evolution of the timbre over time is also restored.

It should also be noted that the analysis and restoration can generally relate to the whole audible frequency band from about 15 Hz to 15 kHz, even if, in practice, it can be limited to 8 kHz. The frequencies of the band from 4 to 8 kHz, cut for standard transmission by the telephone network, are in this case analysed and restored since the corresponding information is transmitted in the form of a remote command from the library 31 which already contains the lines at these frequencies, which avoids any explicit transmission thereof.

It should also be noted that if the analysis can relate to only a limited number of sufficiently characteristic frequency bands in the library 11, 12, the various signals to be restored, in library 32, contain all the lines initially input, i.e. each cover, for example, a single-piece band of 15 Hz to 8 kHz.

As explained at the beginning, the invention can apply in cases not involving transmission, in order, for example to store locally a message which to be restored later, i.e. this would be a tape recording function.

In another embodiment, not illustrated, the primary and secondary code words are associated with facsimile data to form a voice-data message. The message is input via the telephone usually associated with facsimile machines and is restored by the same means at the called facsimile machine. The code words emitted by a circuit such as 28 are inserted in a specific field of the message by a microprocessor managing the facsimile transmission protocol and in the same way are retrieved upon reception to be processed as explained above. It is thus possible to carry out voice annotation of a facsimile message, annotation which is transmitted, for example, as a facsimile header.

CLAIMS

- 1 Process for encoding speech formed from a sequence of acoustic units, in which the units are compared with library references associated with primary code words, the differences between the units and the references are determined, the differences are encoded by secondary code words and pairs of primary and secondary codes are substituted for the units.
- 2 Process according to claim 1, wherein since the comparison relates to the energies of spectra of lines of frequencies, weighting coefficients normed with respect to the energy of the reference lines are determined for the lines, and the said coefficients are integrated into the secondary code word.
- 3 Process according to claim 1, wherein the said difference is determined from a succession of spectra corresponding to a succession of amplitudes of the acoustic unit considered.
- 4 Process according to claim 3, wherein only the amplitudes corresponding to extremes are considered.
- 5 Process according to claim 2, wherein the said difference is determined from a single average spectrum of the acoustic unit considered.
- 6 Process according to claim 2, wherein the frequency comparison is limited to three frequency bands.
- 7 Process according to claim 6, wherein the weighting coefficient of the lines of each band is expressed by a single coefficient.
- 8 Terminal for encoding speech signals, comprising means (26, 27) for inputting a sequence of acoustic units and transmitting it to comparator means (16) arranged to compare successively the acoustic units with references in libraries (11, 12) and thus select therein in each instance a specific primary code word of one of the references, the terminal being characterised by the fact that the comparator means (16) are arranged to determine a



difference between the input acoustic unit considered and the reference corresponding to the code word selected and to transmit this difference to transcoding means (17) provided to supply, in response, a secondary code word corresponding to memory means (18) arranged to associate the respective primary and secondary code words.

- 9 Terminal for decoding speech signals, comprising means (30) for receiving signals representing primary code words of references of acoustic units in a library (32), and decoding means (33) arranged to select certain ones of the references in the library (32) according to the primary code words received and to control a transducer (34) for restoration of speech signals accordingly, the terminal being characterised in that the decoding means (33) are arranged also to decode secondary, correction code words associated with the primary code words, and to correct (38) the selected voice references accordingly.
- 10 Facsimile machine according to claim 8, comprising means to insert the code words into a facsimile message.

ABSTRACT OF THE DISCLOSURE

The process for encoding speech formed from a sequence of acoustic units consists of comparing the units with library references associated with primary code words and determining the differences between the units and the references, the differences are encoded by secondary code words to substitute pairs of primary and secondary codes for the units and, if the comparison relates to the energies of spectra of lines of frequencies, weighting coefficients normed with respect to the energy of the reference lines are determined for the lines, and the coefficients are integrated into the secondary code word.